

Provably Efficient RL under Episode-Wise Safety in Constrained MDPs with Linear Function Approximation

Toshinori Kitamura, Arnob Ghosh, Tadashi Kozuno, Wataru Kumagai,

Kazumi Kasaura, Kenta Hoshino, Yohei Hosoe, Yutaka Matsuo

SINICx

M 松尾研究室
MATSUO LAB THE UNIVERSITY OF TOKYO

 **MOONSHOT**
RESEARCH & DEVELOPMENT PROGRAM

NJIT
New Jersey Institute
of Technology

Contribution

We develop the first **linear CMDP** algorithm that guarantees:

- **episode-wise safe** exploration
- $\tilde{\mathcal{O}}(\sqrt{K})$ regret
- and computational tractability

Table 1: $\tilde{\mathcal{O}}(\sqrt{K})$ regret CMDP algorithms

	paper	Epi.-wise safe?	Comp. Efficient?
Tabular	Yu et al. [47]	Yes	State size dependent
	Ghosh et al. [18]	No	No
Linear	Roknilamouki et al. [36]	Instantaneous	Yes
	Ours	Yes	Yes

Notation: Linear CMDP

Linear CMDP: $(\mathcal{S}, \mathcal{A}, H, P, r, u, b, s_1)$

- Finite state & action spaces: \mathcal{S}, \mathcal{A}
- Horizon: $H \in \mathbb{N}$
- Transition probability: $P_h(s' \mid s, a) \in [0, 1]$
- Reward & utility: $r_h(s, a), u_h(s, a) \in [0, 1]$
- Constraint threshold: $b \in [0, H]$
- Initial state: $s_1 \in \mathcal{S}$

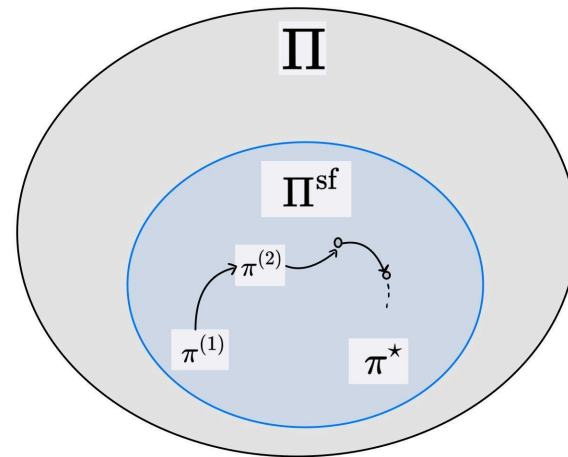
Policy and Value Function

- Value functions: $V_{P,h}^{\pi,r}, V_{P,h}^{\pi,u} : \mathcal{S} \rightarrow \mathbb{R}$, where
$$V_{P,h}^{\pi,r}(s) = \mathbb{E} \left[\sum_{h'=h}^H r_{h'}(s_{h'}, a_{h'}) \mid s_h = s \right]$$
- Action value functions: $Q_{P,h}^{\pi,r}, Q_{P,h}^{\pi,u} : \mathcal{S} \rightarrow \mathbb{R}$, where
$$Q_{P,h}^{\pi,r}(s, a) = \mathbb{E} \left[\sum_{h'=h}^H r_{h'}(s_{h'}, a_{h'}) \mid s_h = s, a_h = a \right]$$

- $P_h(s' \mid s, a) = \boldsymbol{\mu}_h(s')^\top \boldsymbol{\phi}(s, a)$
 - Known feature map $\boldsymbol{\phi} : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}^d$
 - Unknown vectors $\boldsymbol{\mu}_h := (\boldsymbol{\mu}_h^1, \dots, \boldsymbol{\mu}_h^d) \in \mathbb{R}^{S \times d}$
- $r_h(s, a) = (\boldsymbol{\theta}_h^r)^\top \boldsymbol{\phi}(s, a), u_h(s, a) = (\boldsymbol{\theta}_h^u)^\top \boldsymbol{\phi}(s, a)$
 - Known vectors $\boldsymbol{\theta}_r, \boldsymbol{\theta}_u \in \mathbb{R}^d$

Setting: Episode-wise Safe Exploration

- Policies $\pi^{(1)}, \dots, \pi^{(K)} \in \Pi$
- Safe policy set: $\Pi^{\text{sf}} \triangleq \left\{ \pi \mid V_{P,1}^{\pi,u}(s_1) \geq b \right\}$
 - Episode-wise safety: $\pi^{(k)} \in \Pi^{\text{sf}}$ for all k



Assumption (Slater condition): We have access to $\pi^{\text{sf}} \in \Pi^{\text{sf}}$ and $\xi > 0$ such that $V_{P,1}^{\pi^{\text{sf}},u}(s_1) \geq b + \xi$.

Goal: Sublinear regret with episode-wise safety

$$\text{Regret}(K) := \sum_{k=1}^K V_{P,1}^{\pi^*,r}(s_1) - V_{P,1}^{\pi^{(k)},r}(s_1) = o(K) \text{ such that } \pi^{(k)} \in \Pi^{\text{sf}} \quad \forall k \in [1, K]$$

where $\pi^* \in \arg \max_{\pi \in \Pi^{\text{sf}}} V_{P,1}^{\pi,r}(s_1)$ is an optimal safe policy.

Basic Approach: Optimistic-Pessimistic (Opt-Pes) Exploration

Idea (e.g., Yu et al. [47]): Estimate reward value **optimistically** and utility value **pessimistically**:

$$(\text{Opt-Pes}) \quad \pi^{(k)} \in \arg \max_{\pi \in \Pi} \underbrace{\overline{V}_{(k),1}^{\pi,r}(s_1)}_{\geq V_{P,1}^{\pi,r}(s_1)} \text{ such that } \underbrace{V_{(k),1}^{\pi,u}(s_1)}_{\leq V_{P,1}^{\pi,u}(s_1)} \geq b. \quad (1)$$

$\overline{V}_{(k),1}^{\pi,r} := \widehat{V}_{(k),1}^{\pi,r+\beta^{(k)}}$ where $\widehat{V}_{(k),1}^{\pi,r}$ is the value estimate and $\beta^{(k)}$ is the bonus. Similarly, $V_{(k),1}^{\pi,u} := \widehat{V}_{(k),1}^{\pi,u-\beta^{(k)}}$.

- 😊 **Optimism** leads to sublinear regret
- 😊 **Pessimism** ensures $\pi^{(k)} \in \Pi^{\text{sf}}$

Technical contributions

We introduce two main techniques to address these challenges in linear CMDPs:

1. Tractability — The large state space makes optimizing Eq. (1) non-trivial.
2. Feasibility — Pessimistic estimates can make Eq. (1) infeasible.

Technique 1: Efficient Implementation of Opt-Pes Policy

$$(\text{Opt-Pes}) \quad \pi^{(k)} \in \arg \max_{\pi \in \Pi} \overline{V}_{(k),1}^{\pi,r}(s_1) \text{ such that } \underline{V}_{(k),1}^{\pi,u}(s_1) \geq b, \quad (1)$$

Instead of solving Opt-Pes problem, we realize Opt-Pes by "softmax policy" using estimated value functions:

$$\pi_h^{(k),\lambda}(\cdot \mid s) = \text{SoftMax} \left(\overline{Q}_{(k),h}^r(s, \cdot) + \lambda \underline{Q}_{(k),h}^u(s, \cdot) \right)$$

Lemma 6:

- If $\lambda \approx 0$, then $\pi^{(k),\lambda}$ favors **optimistic exploration**
- If $\lambda \gg 0$, then $\pi^{(k),\lambda}$ tries to satisfy the **pessimistic constraint**

$$\overline{Q}_{(k),1}^{r \text{ or } u} := \widehat{Q}_{(k),1}^{\pi,r \pm \beta^{(k)}}$$

We do bisection search in the interval $\lambda \in [0, C_\lambda]$ to find a good λ :

- If $\underline{V}_{(k),1}^{\pi^{(k),\lambda},u}(s_1) \geq b$, then $\pi^{(k),\lambda}$ is safe. Decrease $\lambda \downarrow$
- Otherwise, $\pi^{(k),\lambda}$ is unsafe. Increase $\lambda \uparrow$.

$C_\lambda > 0$ is an upper bound of λ obtained by the Slater condition.

Technique 2: Safe Policy Deployment

$$(\text{Opt-Pes}) \quad \pi^{(k)} \in \arg \max_{\pi \in \Pi} \overline{V}_{(k),1}^{\pi,r}(s_1) \text{ such that } \underline{V}_{(k),1}^{\pi,u}(s_1) \geq b. \quad (1)$$

At the start of training, the bonus may be so large that no policy satisfies the **pessimistic constraint**.

Safe policy deployment technique: Deploy π^{sf} if the policy $\pi^{(k),C_\lambda}$ is unsafe: $\underline{V}_{(k),1}^{\pi^{(k),C_\lambda},u}(s_1) < b$

→ For any k , $\pi^{(k)}$ is ensured to be safe: it is chosen by either

1. π^{sf} that is safe or 2. $\pi^{(k),\lambda}$ that satisfies the pessimistic constraint.

🤔 For sublinear regret, we need to bound the # of π^{sf} deployments...

Theorem 3. With high probability, the safe policy π^{sf} is deployed at most $\tilde{O}(d^3 H^4 \xi^{-2})$ times.

(See Section 2 for an intuitive derivation.)

Algorithm : Optimistic-Pessimistic Softmax Exploration for Linear CMDP

For each $k = 1, 2, \dots, K$, do

1. If $\underline{V}_{(k),1}^{\pi^{(k),C\lambda,u}}(s_1) < b$, all the softmax policies may not be safe. Deploy $\pi^{(k)} = \pi^{\text{sf}}$.
2. Else, find the best λ for the softmax policy:
 - If $\underline{V}_{(k),1}^{\pi^{(k),\lambda,u}}(s_1) \geq b$, λ is large. Decrease $\lambda \downarrow$.
 - Else, increase $\lambda \uparrow$.
 - After sufficient iterations, set $\pi^{(k)} = \pi^{(k),\lambda}$.
3. Sample a trajectory $(s_1^{(k)}, a_1^{(k)}, \dots, s_H^{(k)}, a_H^{(k)})$ by deploying $\pi^{(k)}$.

Main Open Problem: Can we achieve episode-wise safety in linear CMDPs with multiple constraints?

Theorem 4 With high probability, the algorithm achieves

$$\pi^{(k)} \in \Pi^{\text{sf}} \forall k \in \{1, \dots, K\} \quad \text{and} \quad \text{Regret}(K) \leq \tilde{\mathcal{O}}\left(H^4 \xi^{-1} \sqrt{d^5 K}\right)$$